



Calhoun: The NPS Institutional Archive
DSpace Repository

Theses and Dissertations

1. Thesis and Dissertation Collection, all items

1971-06

Investigation of speaker identification based on nasal phonation.

Young, Robert Bryant.

Monterey, California. Naval Postgraduate School

<http://hdl.handle.net/10945/15783>

Downloaded from NPS Archive: Calhoun



Calhoun is the Naval Postgraduate School's public access digital repository for research materials and institutional publications created by the NPS community. Calhoun is named for Professor of Mathematics Guy K. Calhoun, NPS's first appointed -- and published -- scholarly author.

Dudley Knox Library / Naval Postgraduate School
411 Dyer Road / 1 University Circle
Monterey, California USA 93943

<http://www.nps.edu/library>

INVESTIGATION OF SPEAKER IDENTIFICATION
BASED ON NAVAL PHONATION

Robert Bryant Young

United States Naval Postgraduate School



THESIS

INVESTIGATION OF SPEAKER IDENTIFICATION
BASED ON NASAL PHONATION

by

Robert Bryant Young

Thesis Advisor:

J. D. Campbell

June 1971

Approved for public release; distribution unlimited.

Investigation of Speaker Identification

Based on Nasal Phonation

by

Robert Bryant Young
Lieutenant Commander, United States Navy
B.S., Boston College, 1958

Submitted in partial fulfillment of the
requirements for the degree of .

MASTER OF SCIENCE IN ELECTRICAL ENGINEERING

from the

NAVAL POSTGRADUATE SCHOOL
June 1971

ABSTRACT

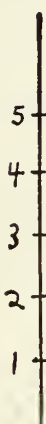
This thesis investigates the possibility of Speaker Identification through the use of Nasal Phonation. Short segments of a restricted set of words from one speaker were sampled, processed, and the resulting vector is used to represent the speaker. Representative vectors were formed for several speakers and correlated with vectors representing individual words from "test" speakers. The magnitude of the correlations of the word vectors with various speaker vectors were used to identify the speaker. This work expands on earlier work done in this field to the extent that it attempts to remove the subjective preparation of data and replace this instead with an objective process of computer mechanization. Some limited success was achieved and, just as important, critical problem areas are noted which, if improved upon as recommended, promise an improved identification capability. Two different word lists fundamental to the identification process were also investigated. Some data was obtained but it was not sufficient to suggest that one word list would be more productive than the other when used as the basis for speaker identification.

Recommendation is made to pursue further research in Speaker Identification using computer programming established during work on this thesis.

during the period of closure, and the vocal cavities remain nearly fixed. Hence, the power spectrum of radiated acoustic energy is essentially steady as indicated by the sound spectrogram of the word "nominal" in Figure 2. Note that the formants show very little movement for the duration of the two nasals. Researchers have commented on the marked extent to which the acoustic properties of nasal consonants vary from speaker to speaker [References 3 and 4]. If these indications are true over a wide population of speakers, it is possible, then, that acoustic radiation produced during the phonation of nasal consonants could provide a strong clue to speaker identity.

Another important feature of nasal consonants with respect to speaker identification is the relative frequency with which they are used in spoken English. According to Tobias [Ref. 5], nasals comprise 11% of the phoneme content of commonly spoken English. Thus, there is sufficient data available in a short sample of speech to provide some degree of speaker identification if our hypothesis is correct. In substantiation of this approach, Glenn and Kleiner [Ref. 1] have reported that in experiments involving a population of ten speakers an average identification accuracy of 97% was obtained. With an experimental population of thirty speakers, identification accuracy was 93%.

FREQUENCY
(KHz)



N O M I N A L

Figure 2. Sonogram of the Word "Nominal".

III. EXPERIMENTAL PROCEDURE

The approach taken in this work is along the lines of Glenn and Kleiner. However, an attempt has been made to use the computer to a greater extent for the processing and manipulation of data and to develop facets of the problem not covered in their experiments.

The basic analysis technique used throughout the experiments described in this thesis is the reduction of a time segment (the near steady state portion during nasal phonation) of selected words into normalized spectral components which form a vector describing a speaker. The same process is again performed on another different set of words spoken by an "unknown" speaker; another vector is formed and is compared with known speaker vectors using the highest value of the cosine of the angle between the compared vectors as the criterion for identification.

A. WORD LIST SELECTION AND RECORDING

Various word lists were used in the different experiments and each will be specifically discussed under the heading of the applicable experiment. However, in each case a word list contained twenty different words beginning with the nasal consonant "n". A given speaker was usually required to record twenty words in slow sequence. The microphone used for this recording was a SURE MODEL 575S. The dynamic range was sufficient for the 1 - 3.5 KHz area of spectrum interest. The recorder used was an Ampex SP-300,

the only one conveniently available. As was later indicated, this recorder caused some problems because of the introduction of high background noise. Since noise bursts could not be tolerated on the tape (because of sensitivity of logic sensing circuits in the digitizing process), considerable care had to be taken during recording. In addition to a rehearsal of the word list, each speaker was cautioned to avoid movements of the microphone, tapping, and other motions, which might introduce such bursts. Also, as it was only possible to use a single recorder output, a procedure was devised whereby the VU meter could be used to judge word placement on the tape during the digitizing process. This procedure required the speaker to hum a tone before and after the word list. This tone could be easily distinguished on the VU meter and greatly facilitated the digitizing process. Several seconds of silence were allowed between words and after the initial tone.

B. DIGITAL CONVERSION

The recorded words were then passed through a 1 - 3.5 KHz Khron-Hite Series 3320 filter and then input to a CI 5000 Analog Computer. Here, after amplification, logic sensing circuits determined the start time of the digitizing process which was performed at 12.5 KHz by an SDS 9300 digital computer.

The start time of the digitizing process is of critical concern because of the need to sample the spoken nasal

consonant during the proper interval -- i.e., when the radiated acoustic energy is most nearly in steady state. As expected, there is a small transient period initially, then near steady state. Fujimuras' analysis [Ref. 3] has shown varying degrees of formant variation toward the tail end of steady state depending on the following vowel. Thus, there appears to be an optimum steady-state window extending for approximately ten milli-seconds during the time of consonant. The logic sensing arrangement allowed for a variable delay to preceed the start of the digitizing process after triggering on the start of a word. Additionally, there was provision of a delay flop to insure that once triggered, the sensing circuit would be immune from further (false) triggers for a predetermined amount of time (experience has shown this time should be about 1.5 seconds).

The number of samples taken was fixed for this processing arrangement at 128 which, at a sampling rate of 12.5 KHz, takes about ten milli-seconds. A ten milli-second delay after triggering was used to insure that the sampling period lay within the steady state optimum sampling window. The choice of sampling frequency also allows a convenient feature (as will be shown later) of having the Fourier coefficients centered in approximately 100 cycle bandwidths.

Thus, the digital output of the conversion for one word list by one speaker consists of twenty blocks of 128 samples each, which are recorded on a seven-track digital

tape as a single file of data. In actual practice, the process was repeated twice again for each file to insure that each word was properly sampled at least once.

In early tests, files were often discarded because it was found that they were not complete. A typical reason for these errors was noise bursts on the voice tape. Though attempts were made to keep this type of error to a minimum, there were still circumstances which required manipulation around these bursts, if data was to be saved. If the burst was distinguishable from the spoken words and occurred prior to the twenty words, it was possible to delay energizing the logic recognition/delay circuits until immediately prior to the first spoken word. The timing was physically difficult to implement in many instances and hence, the requirement to make several runs of the same data to ensure success.

C. TAPE CONVERSION

Since the SDS 9300 did not have the digital storage needed for analysis, use was made of the school's IBM 360. Because the current operating system supports only nine-track tape for FORTRAN input files, it was necessary to convert from the seven-track tape (produced in the digitizing process) to a nine-track tape. To facilitate evaluation of data, a decimal print out of data by block was also provided. Thus, for a given experiment, the entire content of the seven-track tape was converted (included

redundant files); the decimal data was evaluated for completeness, and then, through the use of Job Control Language (JCL), selected files were brought forward and converted for further processing. It should be noted that this editing process allowed mainly a quantitative check on the data; it did not readily permit comparison of this data with the original analog data.

D. POST CONVERSION PROCESSING

Because the energy distribution by frequency was required, the next step in the process was to determine the spectral content of the data blocks. The Fast Fourier Transform Algorithm [Ref. 6] was used to compute the complex Fourier coefficients whose magnitude squared are the energy spectrum. Since a 12.5 KHz sampling rate was used, each value of the energy spectrum represented an incremental bandwidth of 97.66 Hz (very close to the 100 Hz bandwidths that were manually quantized by Glenn and Kleiner). Again, closely paralleling Glenn and Kleiner's work, an approximately 2.5 KHz band of the spectrum was isolated by discarding those values of the energy spectrum below 1025.36 Hz and above 3466.83 Hz. This spectrum band of interest then included twenty-five segments of 97.66 Hz, each of which was considered a component of a twenty-five dimensional vector which represented a particular sampled word. Each of the twenty-word vectors then underwent a normalization transformation according to the formula:

$$v'_i = v_i / \sum_{j=1}^{25} v_j, \quad i = 1, 2, \dots, 25.$$

Additionally another transformation was performed on each vector which was designed to emphasize the major pole and the major zero of the power spectrum. For the vector

$$V' = (v'_1, v'_2, \dots, v'_{25}).$$

Let

$$M = \max \{v'_i\}$$

represent the major pole of the spectrum, and

$$m = \min \{v'_i\}$$

the major zero. Then the transformed vector is given by

$$V^* = (v^*_1, v^*_2, \dots, v^*_{25})$$

where

$$v^*_i = \alpha v'_i - \beta$$

with

$$\alpha = 1/(M-m)$$

$$\beta = m/(M-m).$$

The vectors thus transformed will be called subvectors. At this point, for the purpose of the experiments, the first ten word subvectors were considered to have derived from a known speaker and the second ten word subvectors to have been uttered by an unknown or "test" speaker. Each set of ten subvectors was then arithmetically averaged by component

TABLE IV.

Results of Experiment Three

Test Speaker	Reference Speaker				
	1	2	3	4	5
1	.9164	.1828	.6998	.4863	.7683
2	.3750	.8209	.7028	.6188	.6766
3	.8243	.2121	.6324	.3497	.6580
4	.5117	.7055	.8430	.9299	.8467
5	.9196	.4826	.8894	.7988	.9514

poor correlation is unknown. It was noted in the diary that the speaker's words were close together, and it is possible that the logic circuit did not have sufficient opportunity to stabilize prior to a following word thus introducing false data into the experiment.

The results of Experiment Four are given in Table V. As can be seen, the editing process did improve the match characteristics of the data set. However, the correlation value of speaker five with reference five has now been reduced to a poor third. Speaker three's correlation values have not measurably improved. This might be expected because of the removal of only high correlation subvectors in the editing process.

The results of Experiment Five are shown in Table VI. Again the overall match record is three out of five. However, speaker one was now matched and speaker four unmatched, although it should be noted that speaker four places second. Speaker three is still the least likely match as was the case in the unedited version of the experiment (Experiment Three).

The results of Experiment Six are given in Table VII. There are six clear matches. Of the remaining seven, two miss matches by one, i.e., place second, and one (speaker twelve) placed third. The worst case was the speaker who was mismatched by seven. It is interesting to note that the recording diary had the comment "a little fast" for this speaker's recitation. This means that the speaker

TABLE V.

Results of Experiment Four

Reference Speaker

Test Speaker	1	2	3	4	5
1	.9176	.1175	.4491	.1604	.2906
2	.1854	.9483	.7105	.6391	.7679
3	.9095	.2628	.7076	.4981	.5087
4	.2586	.5276	.8399	.9369	.9526
5	.8627	.3447	.8073	.7110	.7136

TABLE VI.

Results of Experiment Five

Test Speaker	Reference Speaker				
	1	2	3	4	5
1	.9064	.1647	.6169	.5735	.7994
2	.2556	.6771	.7044	.4809	.6503
3	.6148	.1361	.3950	.2903	.6096
4	.3411	.6397	.8069	.7423	.7381
5	.8691	.5194	.8263	.8346	.9542

of a higher percentage of correct identifications, even with a larger population of unidentified speakers.

The most critical point in the digitizing process is the method used to start the digitizing. As was pointed out previously, the start time is sensitive to noise and one must have a check to ensure that it is done properly. Under the present scheme, reliance was placed on a delay flop to start the process. However, what was not taken into consideration was the problem of noise in the tape recorder. It is probable that because of this noise threshold there were instances when the analog voice voltage was beneath the noise for a fraction of time prior to delay timing by the flip-flop circuit. The particular recorder used for the experiments (as mentioned earlier) was noisy and hence, inconsistencies in start of sampling may have resulted.

One way to improve this situation is to use a more noise-free recorder not available for this project. Also, there should be a simultaneous graphical record of the analog voltage from the recorder along with the digital to analog record of the digitized version. For convenience, a milli-second time tick could also be plotted. This would assist the operator on a near real time basis to ensure that data is taken at the desired time. Yet another improvement would be the use of more sophisticated logic to start the digitizing. Such logic might use some type of nasal consonant recognizer to insure that digitizing starts during the steady state.

A second area needing improvement is the editing process. Even with the improvement for start of digitization, there is still the possibility that extraneous recorded sound energy not related to the desired signal will be present in the subvectors. Hence, the need to edit them. In the experiments as described, the first criterion for a "bad" subvector was its poor correlation with the prime vector it helped to generate. These subvectors were discarded mainly on this basis. However, upon closer analysis of the twenty-five component differences between the subvectors and prime vectors, it was found that, although poor correlations could be relied upon in most cases as an indication for rejection, there were other cases where there were irregularities in components which compensated one another and hence, gave a higher correlation value and, therefore, were retained. These latter subvectors should also have been rejected for they were probably sampled for one reason or another during a non steady state condition.

Another area of error overlooked in the editing process was the decision of how many subvectors to eliminate. Three were arbitrarily chosen for each speaker mainly because it was convenient for computer processing and also it was felt that this small number would not drastically affect the true character of each prime vector. However, as was seen in the experiments, when this selection rule was applied to all speaker subvectors, some high correlating

subvectors are removed and hence the data is distorted. The discarding of subvectors should be made solely on the basis of a small magnitude correlation value and, as mentioned above, of component variance.

All speaker data was retained throughout the experiments, and except for the editing described above, no attempt was made to "dress" the results by discarding data which was often misclassified. Speaker three in experiments Two through Five was one that might have been so eliminated. This particular data was suspect because of his fast speech and consequent uncertainty regarding proper digitization. Speaker six in Experiments Six and Seven was another example.

Another interesting fact is that words spoken by speaker two (the author) were always correctly identified (at least prior to the editing process). Though the data base is slim, it suggests that greater care in recording of words spoken by individuals taking part in the experiment would yield measurably improved results.

The best word list to use is still open to question. Data produced from the experiments does not favor either of the two lists studied and in any case is not exhaustive enough for a firm decision.

Lastly, the author would like to point out that, although it was merely a tool to an end, the computer mechanization and debugging of the various programs used for the experiments was formidable. It is hoped that interest in Speaker Identification will continue at the Naval Postgraduate School and that the establishment of the programming

used in this thesis will be used as a stepping stone to further research.

LIST OF REFERENCES

1. Glenn, J. W. and Kleiner, N., "Speaker Identification Based on Nasal Phonation," The Journal of the Acoustical Society of America, v. 43, p. 368-372, February 1968.
2. Flanagan, J. L., Coker, C. H., Rabiner, L. R., Schafer, R. W., and Umeda, N., "Synthetic Voices for Computers," IEEE Spectrum, v. 7, p. 22-45, October 1970.
3. Fujimura, O., "Analysis of Nasal Consonants," The Journal of the Acoustical Society of America, v. 34, p. 1865-1875, December 1962.
4. Dickson, D. R., "An Acoustic Study of Nasality," Journal of Speech and Hearing Research, v. 5, p. 103-111, June 1962.
5. Tobias, V., "Relative Occurrence of Phonemes in American English," The Journal of the Acoustical Society of America, v. 31, p. 631, 1959.
6. Cooley, J. W. and Tukey, J. W., "An Algorithm for the Machine Calculations of Complex Fourier Series," Mathematics of Computations, v. 19, p. 297, April 1965.

INITIAL DISTRIBUTION LIST

	No. Copies
1. Defense Documentation Center Cameron Station Alexandria, Virginia 22314	2
2. Library, Code 0212 Naval Postgraduate School Monterey, California 93940	2
3. Asst Professor J. D. Campbell, Code 52Cb Department of Electrical Engineering Naval Postgraduate School Monterey, California 93940	2
4. LCDR Robert B. Young, USN USNAV SEC GRU ACTY FPO, New York 09513	5

DOCUMENT CONTROL DATA - R & D

(Security classification of title, body of abstract and indexing annotation must be entered when the overall report is classified)

1. ORIGINATING ACTIVITY (Corporate author) Naval Postgraduate School Monterey, California 93940		2a. REPORT SECURITY CLASSIFICATION UNCLASSIFIED	
		2b. GROUP	
3. REPORT TITLE INVESTIGATION OF SPEAKER IDENTIFICATION BASED ON NASAL PHONATION			
4. DESCRIPTIVE NOTES (Type of report and, inclusive dates) Master's Thesis; June 1971			
5. AUTHOR(S) (First name, middle initial, last name) Robert B. Young			
6. REPORT DATE June 1971		7a. TOTAL NO. OF PAGES 42	7b. NO. OF REFS 6
8a. CONTRACT OR GRANT NO.		9a. ORIGINATOR'S REPORT NUMBER(S)	
b. PROJECT NO.			
c.		9b. OTHER REPORT NO(S) (Any other numbers that may be assigned this report)	
d.			
10. DISTRIBUTION STATEMENT Approved for public release; distribution unlimited.			
11. SUPPLEMENTARY NOTES		12. SPONSORING MILITARY ACTIVITY Naval Postgraduate School Monterey, California 93940	
13. ABSTRACT			

This thesis investigates the possibility of Speaker Identification through the use of Nasal Phonation. Short segments of a restricted set of words from one speaker were sampled, processed, and the resulting vector is used to represent the speaker. Representative vectors were formed for several speakers and correlated with vectors representing individual words from "test" speakers. The magnitude of the correlations of the word vectors with various speaker vectors were used to identify the speaker. This work expands on earlier work done in this field to the extent that it attempts to remove the subjective preparation of data and replace this instead with an objective process of computer mechanization. Some limited success was achieved and, just as important, critical problem areas are noted which, if improved upon as recommended, promise an improved identification capability. Two different word lists fundamental to the identification process were also investigated. Some data was obtained but it was not sufficient to suggest that one word list would be more productive than the other when used as the basis for speaker identification.

Recommendation is made to pursue further research in Speaker Identification using computer programming established during work on this thesis.

Nasal Phonation

Speaker Identification

